

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-93286

(43) 公開日 平成7年(1995)4月7日

(51) Int.Cl.⁹

識別記号

庁内整理番号

F I

技術表示箇所

G 0 6 F 17/00

G 0 1 N 33/68

G 0 6 F 17/30

8724-5L

9194-5L

G 0 6 F 15/ 20

15/ 40

D

5 3 0 S

審査請求 有 請求項の数7 O L (全 9 頁)

(21) 出願番号 特願平5-233822

(22) 出願日 平成5年(1993)9月20日

(71) 出願人 000004237

日本電気株式会社

東京都港区芝五丁目7番1号

(72) 発明者 馬見塚 拓

東京都港区芝五丁目7番1号 日本電気株式会社内

(72) 発明者 安倍 直樹

東京都港区芝五丁目7番1号 日本電気株式会社内

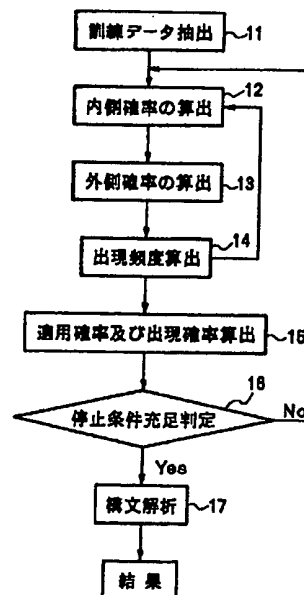
(74) 代理人 弁理士 京本 直樹 (外2名)

(54) 【発明の名称】 タンパク質立体構造予測方法

(57) 【要約】

【目的】 構造未知のタンパク質のアミノ酸配列情報から、その立体構造内に含まれる長距離相互作用に基づく立体構造を高精度で予測する。

【構成】 ステップ11で、構造既知及び未知のタンパク質アミノ酸配列から長距離相互作用を有する立体構造の訓練データを抽出し、ステップ12で、木構造の生成確率における内側確率を算出し、ステップ13で、木構造の生成確率における外側確率を算出し、ステップ14でこれら内側確率及び外側確率から書き換え規則の出現頻度、及び各アミノ酸の出現頻度を算出し、ステップ15で、書き換え規則の適用確率及び規則の末端ノードの各アミノ酸の出現確率を算出し、ステップ16で、訓練データの反復学習の停止条件の充足判定を行い、ステップ17で、学習により得られた書き換え規則を使用して、構文解析により立体構造未知のデータに対して、長距離相互作用を有する立体構造部位の予測を行う。



【特許請求の範囲】

【請求項1】 タンパク質のアミノ酸配列からタンパク質の構造予測を行うための訓練データを抽出するステップと、訓練データからタンパク質の部分的な立体構造に相当する書き換え規則を学習するステップと、学習された書き換え規則を用いて、テストアミノ配列データに対し、立体構造部分の予測を行うステップとからなることを特徴とするタンパク質立体構造予測方法。

【請求項2】 前記訓練データを抽出するステップは、立体構造既知のタンパク質に対し、同じファミリーに属するタンパク質、もしくは一次構造上、一定値以上の相10 同性を有するタンパク質のアミノ酸配列を、アミノ酸配列データベースから抽出することを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項3】 前記書き換え規則を学習するステップにおける書き換え規則が確率的な規則であり、該ステップは、確率的脈自由文法の学習に使用されるインサイド・アウトサイドアルゴリズムと呼ばれる学習方法の本文法への拡張であることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項4】 前記書き換え規則を学習するステップにおける書き換え規則が確率的な規則であり、末端ノードに20種類のアミノ酸もしくは、それらのカテゴリーが割り当てられ、それらの出現確率付き確率的規則であり、該ステップは、確率的文脈自由文法の学習に使用されるインサイド・アウトサイドアルゴリズムと呼ばれる学習方法の本文法への拡張であることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項5】 前記立体構造部分の予測を行うステップは、文脈自由文法の構文解析に使用されるCKYアルゴリズムと呼ばれる構文解析方法の拡張であることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項6】 前記立体構造部分の予測を行うステップにおける書き換え規則が確率的な規則であり、該ステップは、確率的文脈自由文法の構文解析に使用されるCKYアルゴリズムと呼ばれる構文解析方法の本文法への拡張であることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項7】 前記立体構造部分の予測を行うステップにおける書き換え規則が確率的な規則であり、末端ノードに20種類のアミノ酸もしくは、そのカテゴリーが割り10 当てられ、出現確率付き確率的規則であり、該ステップは、確率的文脈自由文法の構文解析に使用されるCKYアルゴリズムと呼ばれる構文解析方法の本文法への拡張であることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、立体構造未知のタンパク質アミノ酸配列から、タンパク質の立体構造を予測す15

る方法に関する。

【0002】

【従来の技術】 タンパク質の立体構造を予測する方法としては、タンパク質全体の立体構造ではなく、その部分的な立体構造である二次構造を予測する方法が一般的である。従来、タンパク質二次構造予測問題は、タンパク質の一次構造の各残基（以下、予測対象となる残基を中心残基と呼ぶ）が、 α ヘリックス、 β シート、それ以外という3種類の二次構造のいずれに対応するかを予測する問題として扱われてきた。従来技術によるタンパク質の二次構造を予測する方法として、例えば、1974年発行の米国の雑誌「バイオケミストリー」(Biochemistry)の第23巻222-245頁記載のチヨウ(Chou)とファスマン(Fasman)による論文「プレディクション オブ プロテイン コンホメーション」(Prediction of protein conformation)（以下、CF法と略す）、1978年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」(Journal of Molecular Biology)の第120巻97-120頁掲載のガルニエ(Garnier)らによる論文「アナリシス オブ ザ アクキュレシー アンド インプリケーションズ オブ シンプル メソッド フォー プレディクティング ザ セコンダリー ストラクチャー オブ グロブラー プロテインズ」(Analysis of the accuracy and implications of simple method for predicting the secondary structure of globular proteins)（以下、GOR法と略す）、1987年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」(Journal of Molecular Biology)の第198巻425-443頁掲載のギブラト(Gibrat)らによる論文「ファザー デベロプメンツ オブ プロテイン セコンダリー ストラクチャー プレディクション ユージング インホメーション セオリー：ニュー パラメータズ アンド コンシダレーション オブ レジデューペアズ」(Further developments of protein secondary structure prediction using information theory: New parameters and consideration of residue pairs)（以下、GGR法と略す）、1988年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」(Journal of Molecular Biology)の第202巻865-884頁記載のキャン(Qian)らによる論文「プレディクティング ザ セコンダリー ストラクチャー オブ グロブラー プロテインズ ユ20

ージング ニューラル ネットワーク モデルズ」(Predicting the secondary structure of globular proteins using neural network models) (以下、QS法と略す) などがある。

【0003】CF法は、タンパク質構造のデータベースから各二次構造におけるアミノ酸の統計的な出源頻度を求め、この頻度表を使用し、経験的な規則に基づく予測を行っている。また、GOR法は、中心残基の二次構造に対して、その残基から数残基離れた残基により独立にもたらされる情報量の和を計算し、その相対値から予測を行い、GGR法は、中心残基の二次構造に対して、その残基及びその残基から数残基離れた残基によりもたらされる情報量の和から予測を行っている。さらにQS法は、3層のフィードフォワード型のネットワークを使用し、中心残基の前後8残基を含む配列を入力とし、二次構造に対する中心残基及び周辺残基からの寄与をニューラルネットワークを用いて抽出することにより予測を行っている。

【0004】

【発明が解決しようとする課題】タンパク質の立体構造の中には、その部分構造においても、数残基から数十残基以上も離れたアミノ酸残基同士の相互作用(以下、遠距離相互作用と呼ぶ)により構成されているものが数多くある。例えば、主要な二次構造の一つである β シートも、その一つである。しかし、従来、タンパク質立体構造予測、例えば、 β シートの予測を含む二次構造予測手法などにおいて、そのような遠距離相互作用を保持している残基同士の間に存在するアミノ酸残基を無視し、遠距離相互作用を自動的に抽出する方法や、さらに、その遠距離相互作用を規則として未知データに対する予測を行う方法は皆無であり、そういった手法は確立されていなかった。

【0005】

【課題を解決するための手段】本発明のタンパク質立体構造予測方法は、タンパク質のアミノ酸配列からタンパク質の構造予測を行うための訓練データを抽出するステップと、訓練データからタンパク質の部分的な立体構造に相当する書き換え規則を学習するステップと、学習された書き換え規則を用いて、テストアミノ酸配列データに対し、立体構造部分の予測を行うステップとからなることを特徴とする。

【0006】又、前記訓練データを抽出するステップは、立体構造既知のタンパク質に対し、同じファミリーに属するタンパク質、もしくは一次構造上、一定値以上の相同性を有するタンパク質のアミノ酸配列を、アミノ酸配列データベースから抽出することを特徴とする。

【0007】又、前記書き換え規則を学習するステップにおける書き換え規則が確率的な規則であり、該ステップは、確率的文脈自由文法の学習に使用されるインサイ

ド・アウトサイドアルゴリズムと呼ばれる学習方法の本文法への拡張であることを特徴とする。

【0008】又、前記書き換え規則を学習するステップにおける書き換え規則が確率的な規則であり、末端ノードに20種類のアミノ酸もしくは、それらのカテゴリーが割り当てられ、それらの出現確率付き確率的規則であり、該ステップは、確率的文脈自由文法の学習に使用されるインサイド・アウトサイドアルゴリズムと呼ばれる学習方法の本文法への拡張であることを特徴とする。

【0009】又、前記立体構造部分の予測を行うステップは、文脈自由文法の構文解析に使用されるCKYアルゴリズムと呼ばれる構文解析方法の拡張であることを特徴とする。

【0010】又、前記立体構造部分の予測を行うステップにおける書き換え規則が確率的な規則であり、該ステップは、確率的文脈自由文法の構文解析に使用されるCKYアルゴリズムと呼ばれる構文解析方法の本文法への拡張であることを特徴とする。

【0011】又、前記立体構造部分の予測を行うステップにおける書き換え規則が確率的な規則であり、末端ノードに20種類のアミノ酸もしくは、そのカテゴリーが割り当てられ、出現確率付き確率的規則であり、該ステップは、確率的文脈自由文法の構文解析に使用されるCKYアルゴリズムと呼ばれる構文解析方法の本文法への拡張であることを特徴とする。

【0012】

【実施例】次に、本発明について図面を参照して詳細に説明する。本実施例では、対象とする遠距離相互作用からなるタンパク質立体構造として β シート領域を扱うものとする。

【0013】図1は、本発明において遠距離相互作用を保持した立体構造の規則に相当する書き換え規則の例である。一般に、書き換え規則は、非終端記号及び終端記号がラベル付けされたノードを有する木であり、固有の適用確率を持つ。書き換え規則の書き換え動作は、あらかじめ設定した初期木に対して、書き換え規則を有限回適用し、終端記号のみからなる木構造を構成することで達成される。

【0014】ここで、木構造とは、単一ノード、もしくは複数の木構造を左右に順序づけられた子供として持つノードからなる構造として再帰的に定義される。また、木構造のルートノードとは、そのノードを子供とするノードが存在しない木構造の中の唯一のノードを指す。さらに、ノードのランクとは、ノードの下の子供の数を指す。例えば、図1はランク1の書き換え規則である。

【0015】具体的に、生成確率 P_A を持つ木構造Aに含まれるランクrの非終端記号によりラベル付けされたノードTを、適用確率 P_B を有する同ランクの木構造Bによって書き換えるということは、A中のTをBによって置き換え、Tのr個の子供を各々B中のr個の空白ノ

ードの位置に、左から順番に連結し、新たに導出された木構造とし、さらに、 P_A と P_B との積をとることにより、新しい生成確率を計算することを指す。

【0016】ここで、終端記号は、20種類のアミノ酸、もしくは、アミノ酸をその化学的性質に基づいて分類したいくつかのグループに対応し、導出された木構造の末端に現れるそれらアミノ酸もしくはグループの集合はアミノ酸配列に相当する。与えられたアミノ酸配列に対し、適当な書き換え規則を使用し、そのアミノ酸配列が木構造の末端に出現するように書き換え動作を行うことにより、その配列上のどの部分が各書き換え規則により生成されたかを調べることを「構文解析 (parsing)」という。

【0017】また、タンパク質アミノ酸配列の各残基位置においては、1つのアミノ酸種類に固定されず、他のアミノ酸に置き換えられても、その立体構造及び機能を保持していることが多い。そのため、一般的な書き換え規則をそのまま使用した場合、全体では類似していながら、末端文字のみが異なる木構造が数多く出現する。そこで、あらかじめ、書き換え規則を、その末端にラベル付けされている一つの終端記号を20種類のアミノ酸と各々に付随した20の出現確率とにそれぞれ置き換え、出現確率を加味した書き換え規則（以下、出現確率付き書き換え規則）にしておいた方が、書き換え規則の数が減り、計算上都合が良く、また、規則も見やすくなる。この場合、生産確率 P_A を持つ木構造Aに含まれるランク r の非終端記号によりラベル付けされたノードTを、適用確率 P_B を有する同ランクの出現確率付き確率規則Bによって書き換えた場合、 P_A と P_B とBの各末端ノードの対応するアミノ酸に付随した出現確率の積をとることにより、新しい生成確率を計算する。出現確率付き書き換え規則の例を図2に示す。

【0018】以下、ランク1の出現確率付き書き換え規則において、1種類の非終端記号をルートノードに使用する場合について、 β シート領域の規則の学習方法、及び規則の適用による予測方法について具体的に説明する。

【0019】図3は、本発明のタンパク質立体構造予測方法の実施例を説明するフローチャートである。

【0020】ステップ11では、まず、 β シート領域既知のタンパク質をタンパク質立体構造データベースから抽出する。さらに、該タンパク質に対し、同じタンパク質あるいは同じファミリーに属するタンパク質、もしくは、一次構造上、一定の割合で相同性を保持しているタンパク質のアミノ酸配列をタンパク質アミノ酸配列データベースから抽出し、学習に使用するデータとする。配列データベースから抽出したタンパク質においては、 β シート領域が未知でも構わないとする。

【0021】例えば、イミノグロブリンというタンパク質において、ヒトのそのタンパク質のアミノ酸配列上の

どこに β シート領域が存在するかは物理化学的な実験から明らかになっている。このヒト・イミノグロブリンに対して、ヒト以外の種類、例えば、チンパンジー、イヌ等のイミノグロブリンのアミノ酸配列、あるいは、アラメントにより一定の割合以上の相同性を有するアミノ酸配列を学習データとする。

【0022】ステップ12、13、14、及びステップ15は、ステップ11で得られた学習データを使用し、あらかじめ設定した構造を有する書き換え規則の適用確率及び書き換え規則の末端ノードにおける各アミノ酸の出現確率を学習するステップである。

【0023】ステップ12は内側確率の算出を行うステップ、ステップ13は外側確率の算出を行うステップ、ステップ14は内側確率、外側確率から、書き換え規則の適用確率及びその末端ノードにおける各アミノ酸の出現確率の計算に必要な、書き換え規則の出現頻度及び末端ノードの各アミノ酸の出現頻度を計算するステップである。ステップ15は、算出された出現頻度を使用し、書き換え規則の適用確率及びその末端ノードの各アミノ酸の出現確率を計算する。

【0024】ステップ12の内側確率の算出過程を、図4に示したフローチャートを用いて説明する。内側確率の算出には、4次元のテーブル $In[i, j, k, l]$ を用意する。テーブル In の次元は書き換え規則のランクに依存し、書き換え規則のランク r に対し、 $2(r+1)$ 次元のテーブルが必要である。テーブル $In[i, j, k, l]$ においては、有限回の書き換え規則の適用により生成された木構造の末端に出現するアミノ酸配列が、与えられた訓練データの i 番目から j 番目までの残基位置、 k 番目から l 番目までの残基位置に相当しており、それらすべての木構造が生成された生成確率の和を示す。訓練データの各配列が与えられる前に、それぞれ $In[i, j, k, l] = 0.0$ に設定しておく。

【0025】訓練データの配列が与えられた場合に、ステップ21において、テーブル $In[i, j, k, l]$ を初期化する。初期化は、訓練データに対し、書き換え規則を1回だけ適用し、得られた木構造の末端ノードと訓練データの部分配列とを対応させることを指す。この動作により、生成された木構造の末端に出現するアミノ酸配列が、与えられたアミノ酸配列の部分配列に相当し、対応するテーブル $In[i, j, k, l]$ に生成確率を格納する。もし、他の書き換え規則の適用により生成された部分配列に相当する。アミノ酸配列の位置が、同様に i, j, k, l であれば、この書き換え規則の生成確率を、テーブル $In[i, j, k, l]$ に加算する。

【0026】例えば、長さが4以上のアミノ酸配列に、図2に示す書き換え規則を適用した場合に生成された部分配列に相当するアミノ酸配列の位置としては、例えば、 $i=1, j=2, k=3, l=4$ が考えられる。この

10

20

30

40

50

時、書き換え規則の各末端ノードに出現するアミノ酸の出現確率及び書き換え規則の適用確率のすべての積をとったものがテーブルIn [1, 2, 3, 4] の値となる。

【0027】以上のように、与えられた訓練データにおいて、書き換え規則の適用により生成された部分配列に相応する、取り得るすべてのアミノ酸配列の位置のテーブルの初期化を行う。

【0028】次に、与えられた訓練データに対して取り得るすべてのテーブルInに格納する値の計算を行う。

【0029】まず、ステップ22において、訓練データのアミノ酸配列の長さがNであれば、 $i=N$, $j=i$, $k=N$, $l=k$ と設定する。さらに、ステップ23において、 i, j, k, l の値を動かしながら、各 i, j, k, l の値において、ステップ24の動作を行う。

【0030】ステップ23における i, j, k, l の動作を説明する。 i を1になる($i=1$)まで1ずつ減らし($i=i-1$)、各 i において j をNになる($j=N$)まで1ずつ増やし($j=j+1$)、各 j において k を j になる($k=j$)まで1ずつ減らし($k=k-1$)、各 k において l がNになる($l=N$)まで、 l の値を1ずつ増や($l=l+1$)していく。

【0031】ステップ24においては、各 i, j, k, l で、すべての書き換え規則の末端ノード数を調べ、テーブルIn [i, j, k, l]への書き換え規則の適用により得られた木構造に相当するテーブルInの値が0.0ではない場合にのみ、生成確率の計算を行ないテーブルIn [i, j, k, l]に格納する。もし、複数の書き換え規則により、テーブルIn [i, j, k, l]の生成確率が算出されれば、それらの和をテーブルIn [i, j, k, l]に格納する。

【0032】例えば、 $i=3, j=7, k=8, l=10$ であり、ある書き換え規則の末端ノード数がそれぞれ1, 2, 0, 1であれば、テーブルIn [4, 5, 8, 9]に0.0でない値が格納されている時、訓練データのアミノ酸配列のそれぞれ3, 6, 7, 10番目のアミノ酸に対応する末端ノードでの各アミノ酸の出現確率と書き換え規則の適用確率とテーブルIn [4, 5, 8, 9]の積を計算し、その値を生成確率としてテーブルIn [3, 7, 8, 10]に加算する。

【0033】ステップ25では、ステップ24の動作の終了の判断を行う。ステップ23により i, j, k, l の値が $i=1, j=N, k=N, l=N$ となった場合、ステップ24の動作後、内側確率の計算を終了する。

【0034】以上の動作により、取り得るすべてのIn [i, j, k, l]が計算でき、ステップ12の内側確率の算出を終了する。

【0035】次に、ステップ13での外側確率の算出過程を図5に示したフローチャートを用いて説明する。

【0036】内側確率の算出と同様に、外側確率の算出

においても、4次元のテーブルOut [i, j, k, l]を使用する。テーブルOutの次元は、Inと同様に書き換え規則のランクに依存する。テーブルOut

[i, j, k, l]は、有限回の書き換え規則の適用により生成された木構造の末端ノードに出現するアミノ酸配列が、与えられたN残基からなる訓練データの1番目から i 番目までの残基位置、 j 番目から k 番目までの残基位置、 l 番目からN番目までの残基位置に相当しており、生成された木構造の生成確率の和を示す。各テーブルOut [i, j, k, l]は、訓練データの各配列が与えられる前に、0.0に設定しておく。

【0037】訓練データの配列が与えられた場合に、ステップ31において、テーブルOut [i, j, k, l]を初期化する。初期化は、初期木に対して書き換え規則を1回のみ適用することを指す。この動作により、生成された木構造の末端に出現するアミノ酸配列が、与えられたアミノ酸配列の部分配列に相当し、対応するテーブルOut [i, j, k, l]に生成確率を格納する。もし、他の書き換え規則の適用により生成された部分配列の訓練データ上の位置が、同様に i, j, k, l であれば、この書き換え規則各末端ノードに出現するアミノ酸の出現確率及び書き換え規則の適用確率のすべての積を、テーブルOut [i, j, k, l]に加算する。

【0038】例えば、長さが10の訓練データに、図2に示す書き換え規則の適用により生成された部分配列に相当するアミノ酸配列の位置は、 $i=1, l=10$ でなければならず、さらに j, k に関しては、 $k=j+1$ を満たす7通りが考えられる。例えば、 $i=1, j=5, k=6, l=10$ であれば、各位置のアミノ酸が書き換え規則の各末端ノードに出現する出現確率及び書き換え規則の適用確率のすべての積を算出したものをテーブルOut [1, 5, 6, 10]に加算する。

【0039】以上のように、与えられた訓練データにおいて、書き換え規則の適用により生成された部分配列に相当する、取り得るすべてのアミノ酸配列の位置に対応するテーブルOutの初期化を行う。

【0040】次に、与えられた訓練データに対して取り得るすべてのテーブルOutに格納する値の計算を行う。

【0041】まず、ステップ32において、訓練データのアミノ酸配列の長さがNであれば、 $i=1, j=N, k=j, l=N$ と設定する。さらに、ステップ33において、 i, j, k, l の値を動かしながら、各 i, j, k, l において、ステップ34の動作を行う。

【0042】ステップ33における i, j, k, l の動作を説明する。 i をNになる($i=N$)まで1ずつ増やし($i=N$)、各 i において j を i になる($j=i$)まで1ずつ減らし($j=j-1$)、各 j において k をNになる($k=N$)まで1ずつ増やし($k=k+1$)、各 k

において1がkになる ($1=k$) まで、1の値を1ずつ減ら ($1=1-1$) していく。

【0043】ステップ34においては、各書き換え規則の末端ノード数を調べ、その書き換え規則を適用した場合に、生成された木構造がテーブルOut [$i, j, k, 1$] に対応するような木構造が存在する場合のみ、生成確率の計算を行ないテーブルOut [$i, j, k, 1$] に格納する。もし複数の書き換え規則により、テーブルOut [$i, j, k, 1$] の生成確率が算出されれば、それらの和をテーブルOut [$i, j, k, 1$] に格納する。

【0044】例えば、 $i=3, j=4, k=7, l=9$ であり、書き換え規則の末端ノード数がそれぞれ1, 2, 0, 1であれば、テーブルOut [$2, 6, 7, 10$] に値が格納されている時、訓練データのアミノ酸配列のそれぞれ3, 4, 5, 9番目のアミノ酸に対応する末端ノードの各アミノ酸の出現確率と書き換え規則の適用確率とOut [$2, 6, 7, 10$] の積を計算し、その値を生成確率としてテーブルOut [$3, 4, 7, 9$] に加算する。

【0045】ステップ35では、ステップ34の動作の終了の判断を行う。ステップ33により、 i, j, k, l の値が $i=N, j=i, k=N, l=k$ となった場合、ステップ34の動作後、外側確率の計算を終了する。

【0046】以上の動作により、取り得るすべてのOut [i, j, k, l] が計算でき、ステップ13の外側確率の算出を終了する。

【0047】次に、ステップ14において、内側確率、外側確率から、書き換え規則の出現頻度、及びその末端ノードの各アミノ酸の出現頻度を計算する。

【0048】ステップ14の出現頻度の算出過程を図6に示したフローチャートを用いて説明する。まず、出現頻度を格納する4次元のテーブルPd [m, n, p, q] を用意する。テーブルPdの次元は、In及びOutとは異なり、書き換え規則のランクに依存しない。20種類のアミノ酸を1から20までのアミノ酸番号に置き換えた場合に、テーブルPd [m, n, p, q] の添え字 m, n, p, q は、 m 番目の書き換え規則において、非終端記号の n 番目のノード位置の p 番目の末端ノードに出現する q 番目のアミノ酸を示す。各テーブルPd [m, n, p, q] の値は、訓練データの最初の配列が与えられた時にのみあらかじめ0.0に設定しておく。

【0049】ステップ41においては、 $i=1, j=N, k=j, l=N$ と設定する。さらに、ステップ42において i, j, k, l の値を動かしながら各 i, j, k, l において、ステップ43の動作を行う。

【0050】ステップ42での i, j, k, l の動作を説明する。 i を N になる ($i=N$) まで1ずつ増やし

($i=N$)、各 i において j を i になる ($j=i$) まで1ずつ減らし ($j=j-1$)、各 j において k を N になる ($k=N$) まで1ずつ増やし ($k=k+1$)、各 k において l が k になる ($l=k$) まで、1の値を1ずつ減ら ($l=1-1$) していく。

【0051】ステップ43においては、各書き換え規則それぞれに対し、テーブルOut [i, j, k, l] に対応する木構造に書き換え規則を適用した場合に、対応する添え字を要素とするテーブルInが0.0ではない値を有して存在するかどうかをチェックする。存在していれば、そのテーブルInとOut [i, j, k, l] との間を埋める部分配列に対応する書き換え規則の各末端ノードの各アミノ酸の出現確率、書き換え規則の適用確率、テーブルOut [i, j, k, l]、Inの積を計算し、テーブルPdに加算する。

【0052】具体的に、 $i=2, j=6, k=7, l=12$ の時、図2の構造をした適用確率Pを有する書き換え規則を1番として適用する場合を考える。図2の書き換え規則においては、4個のノード位置の一つずつノードが存在するので、各々訓練データの3, 5, 8, 11番目の残基位置のアミノ酸に対応する。また、この位置のアミノ酸はそれぞれ、アミノ酸番号により、2, 15, 18, 7番であり、書き換え規則の対応する各ノードのアミノ酸の出現確率は、それぞれ p_1, p_2, p_3, p_4 であるとする。この時、Pd [$1, 1, 1, 2$], Pd [$1, 2, 1, 15$], Pd [$1, 3, 1, 18$], Pd [$1, 4, 1, 7$] それぞれに、 $P \times p_1 \times p_2 \times p_3 \times p_4 \times \text{Out} [2, 6, 7, 12] \times \text{In} [4, 4, 9, 10]$ が加算される。

【0053】ここで、テーブルInの添え字は、テーブルOutの添え字と書き換え規則のノードの単純な差ではないことに注意する。

【0054】ステップ44では、ステップ43の動作の終了の判断を行う。ステップ42により、 i, j, k, l の値が $i=N, j=i, k=N, l=k$ となった場合、ステップ43の動作後、テーブルPdの算出を終了する。

【0055】訓練データの各配列に対して、ステップ12、13、14を繰り返し、書き換え規則及び規則の各末端ノード位置での各アミノ酸の出現頻度を算出する。

【0056】ステップ15では、算出を行ったテーブルPdから、書き換え規則の適用確率、及び各末端ノードの各アミノ酸の出現確率を計算する。

【0057】M番目の書き換え規則の適用確率は、

【0058】

【数1】

$$\frac{\sum_m \sum_n \sum_p Pd[m, n, p, q]}{\sum_m \sum_n \sum_p Pd[m, n, p, q]}$$

【0059】により計算される。また、M番目の書き換

え規則のN番目のノード位置のP番目のノードにQ番の
アミノ酸が出現する出現確率は、

【0060】

【数2】

$$\frac{Pd[M, N, P, Q]}{\sum_q Pd[M, N, P, q]}$$

【0061】により計算される。

【0062】以上により、訓練データセットに対して、あらかじめ設定した書き換え規則の適用確率及び各書き換え規則の末端ノード位置にアミノ酸が出現する出現確率が算出される。

【0063】ステップ12、13、14、15までをあらかじめ設定した回数、もしくは、あらかじめ設定した停止条件を満たすまで繰り返す。ステップ16は、停止条件が満たされているかどうかをチェックする。例えば、停止条件としては、「いずれの書き換え規則の適用確率やアミノ酸の出現確率の値も変化が0.01未満である」などが採用できる。

【0064】ステップ17では、与えられたテストアミノ酸配列に対し、書き換え規則の適用確率及び規則の各末端ノードのアミノ酸の出現確率から、構文解析により、どの書き換え規則の末端ノードがテスト配列の部分領域に対応するかを調べる。すなわち、構文解析により、遠距離相互作用を保持した書き換え規則がテストアミノ酸配列のどの部分に現れたかを検出し、テストアミノ酸配列内の遠距離相互作用を保持している部分を予測する。

【0065】ステップ17の構文解析方式を、図7に示したフローチャートを用いて説明する。まず、4次元のテーブルPar[i, j, k, l]を用意する。テーブルParの次元は書き換え規則のランクに依存し、書き換え規則のランクrに対し、2(r+1)次元のテーブルを用意する。テーブルPar[i, j, k, l]は、書き換え規則の適用により生成された木構造の末端ノードに出現するアミノ酸配列が、与えられたN残基からなるテストデータの1番目からi番目までの残基位置、j番目からk番目までの残基位置、1番目からN番目までの残基位置に相当しており、生成された木構造の生成確率の最大値を示す。テストデータが与えられる前に、各テーブルPar[i, j, k, l] = 0.0と設定しておく。

【0066】テストアミノ酸配列が与えられた場合に、ステップ51において、テーブルPar[i, j, k, l]を初期化する。初期化は、初期木に対して書き換え規則を1回のみ適用することを指す。この動作により、生成された木構造の末端ノードのアミノ酸配列が、与えられたアミノ酸配列の部分配列に相当し、対応するテーブルPar[i, j, k, l]に生成確率を格納する。もし、他の書き換え規則の適用により生成された部分配列

のテストデータ上の位置が、同様にi, j, k, lであれば、この書き換え規則の各末端ノードに出現するアミノ酸の出現確率及び書き換え規則の適用確率の積を、Par[i, j, k, l]と比較し、大きい方をPar[i, j, k, l]の値とする。

【0067】次に、テストデータに対して、取り得るすべてのテーブルOut[i, j, k, l]に格納する値の計算を行う。

【0068】まず、ステップ52において、訓練データのアミノ酸配列の長さがNであれば、i=1, j=N, k=j, l=Nと設定する。さらに、ステップ53において、i, j, k, lの値を動かしながら、各i, j, k, lにおいて、ステップ54の動作を行う。

【0069】ステップ53におけるi, j, k, lの動作を説明する。iをNになる(i=N)まで1ずつ増やし(i=N)、各iにおいてjをiになる(j=i)まで1ずつ減らし(j=j-1)、各jにおいてkをNになる(k=N)まで1ずつ増やし(k=k+1)、各kにおいてlがkになる(l=k)まで、lの値を1ずつ減らし(l=l-1)していく。

【0070】ステップ54においては、各書き換え規則の末端ノード数を調べ、その書き換え規則を適用した場合に、生成された木構造がテーブルPar[i, j, k, l]に対応するような木構造が存在する場合のみ、生成確率の計算を行ないテーブルPar[i, j, k, l]に格納する。もし、複数の書き換え規則により、テーブルPar[i, j, k, l]の値が算出されれば、それらの中で最大の値をテーブルPar[i, j, k, l]に格納する。

【0071】ステップ55では、ステップ54の動作の終了の判断を行う。ステップ33により、i, j, k, lの値がi=N, j=i, k=N, l=kとなった場合、ステップ54の動作後、Parの計算を終了する。

【0072】ステップ56では、得られた木構造の末端ノードからなるアミノ酸配列が、与えられたテストデータに対応しているテーブルParの中で、最大のParを選出する。

【0073】ステップ57では、最大のParにおいて、それを算出するために、どの書き換え規則が使われてきたかをチェックする。このチェックは、ステップ54におけるParの算出において、軌跡を記憶しておくことによっても達成される。チェックの結果、例えば、テスト配列に対し、βシート領域に相当する書き換え規則が適用されていれば、適応により生成された部分配列は、βシート領域とアミノ酸配列レベルで非常に近い性質を保持しており、βシート領域である可能性が高いと予測する。

【0074】

【発明の効果】立体構造の既知のタンパク質のアミノ酸配列情報から、立体構造未知のタンパク質の遠距離相互

作用に由来する立体構造を従来技術に対して高い精度で予測可能である。すなわち、本手法により、遠距離に位置するアミノ酸残基同士の相互作用を、中間領域を介せず抽出可能であり、既存手法による局所領域からの予測では誤って予測されていたような領域を相互作用の有無という観点からより实际的に予測することが可能になった。また、本手法は β シートなどの部分的な立体構造として著名な部分のみならず、一次構造上離れた残基同士の相互作用により構成されている活性部位などの機能部位の特徴配列を抽出し、規則として予測に使用することが可能である。

【図面の簡単な説明】

【図1】本発明で使用する書き換え規則の一例を示す模式図

【図2】本発明で使用する書き換え規則の一例を示す模式図

【図3】本発明のタンパク質立体構造予測の一実施例を示すフローチャート

【図4】本発明の学習方式の一部である内側確率の算出方法の一実施例を示すフローチャート

【図5】本発明の学習方式の一部である外側確率の算出方法の一実施例を示すフローチャート

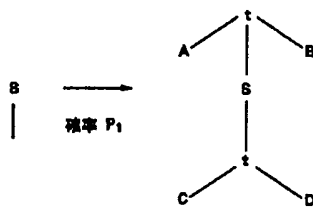
【図6】本発明の学習方式の一部である書き換え規則の出現頻度、及び規則の末端ノードの各アミノ酸の出現頻度の算出方法の一実施例を示すフローチャート

【図7】本発明の構文解析方式の一実施例を示すフローチャート

【符号の説明】

- 11 訓練データ抽出
- 12 内側確率の算出
- 13 外側確率の算出
- 14 出現頻度算出
- 15 適用確率及び出現確率算出
- 16 停止条件充足判定
- 17 構文解析
- 21 内側確率の初期化
- 22 初期添え字の設定
- 23 添え字の更新
- 24 内側確率の算出
- 25 停止条件充足判定
- 31 外側確率の初期化
- 32 初期添え字の設定
- 33 添え字の更新
- 34 外側確率の算出
- 35 停止条件充足判定
- 41 初期添え字の設定
- 42 添え字の更新
- 43 出現頻度の算出
- 44 停止条件充足判定
- 51 最大生成確率の初期化
- 52 初期添え字の設定
- 53 添え字の更新
- 54 最大生成確率の算出
- 55 停止条件充足判定
- 56 最大生成確率の算出
- 57 最大生成確率時の軌跡の検出

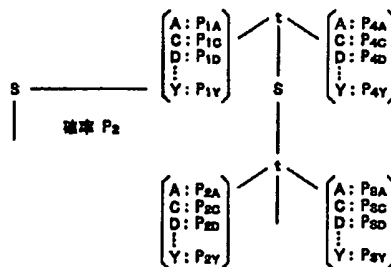
【図1】



S: 非終端記号
A, B, C, D: 終端記号 (アミノ酸)
t: ノード

書き換え規則の例

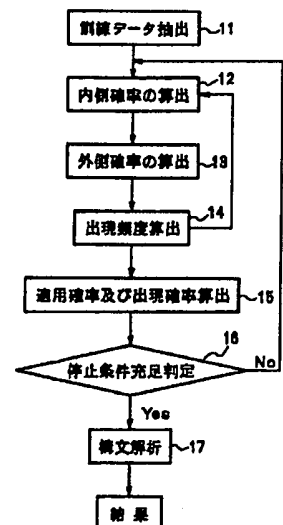
【図2】



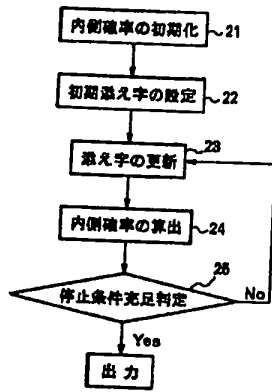
S: 非終端記号
A, C, D, ... Y: 終端記号 (アミノ酸)
t: ノード

出現確率つきの書き換え規則の例

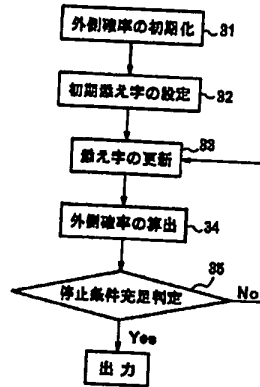
【図3】



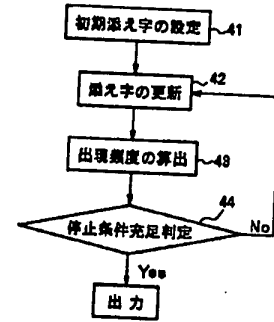
【図4】



【図5】



【図6】



【図7】

